

**УДК 004.4  
ББК 32.372  
Х20**

- Харенслак Б., де Руйтер Дж.**
- X20 Apache Airflow и конвейеры обработки данных / пер. с англ. Д. А. Беликова. – М.: ДМК Пресс, 2022. – 502 с.: ил.

**ISBN 978-5-97060-970-5**

Конвейеры обработки данных управляют потоком данных с момента их первоначального сбора до консолидации, очистки, анализа, визуализации и многого другого. Эта книга научит вас создавать и сопровождать эффективные конвейеры обработки данных с использованием платформы Apache Airflow.

Те, кто мало знаком с Airflow, получат базовое представление о принципах работы этой платформы в I части книги. Далее обсуждаются такие темы, как создание собственных компонентов, тестирование, передовые практики и развертывание, – эти главы можно читать в произвольном порядке в зависимости от конкретных потребностей читателя.

Издание предназначено для специалистов по DevOps, обработке и хранению данных, машинному обучению, а также системных администраторов с навыками программирования на Python.

УДК 004.4  
ББК 32.372

Original English language edition published by Manning Publications USA. Russian-language edition copyright © 2021 by DMK Press. All rights reserved.

Все права защищены. Любая часть этой книги не может быть воспроизведена в какой бы то ни было форме и какими бы то ни было средствами без письменного разрешения владельцев авторских прав.

# Оглавление

---

Часть I ■ ПРИСТУПАЕМ К РАБОТЕ .....	25
1 ■ Знакомство с Apache Airflow.....	27
2 ■ Анатомия ОАГ .....	46
3 ■ Планирование в Airflow.....	67
4 ■ Создание шаблонов задач с использованием контекста Airflow .....	89
5 ■ Определение зависимостей между задачами.....	114
Часть II ■ ЗА ПРЕДЕЛАМИ ОСНОВ .....	144
6 ■ Запуск рабочих процессов .....	146
7 ■ Обмен данными с внешними системами .....	166
8 ■ Создание пользовательских компонентов .....	190
9 ■ Тестирование .....	222
10 ■ Запуск задач в контейнерах .....	259
Часть III ■ AIRFLOW НА ПРАКТИКЕ .....	294
11 ■ Лучшие практики.....	295
12 ■ Эксплуатация Airflow в промышленном окружении .....	324
13 ■ Безопасность в Airflow .....	369
14 ■ Проект: поиск самого быстрого способа передвижения по Нью-Йорку .....	393
Часть IV ■ ОБЛАКО .....	415
15 ■ Airflow и облако .....	417
16 ■ Airflow и AWS .....	426
17 ■ Airflow и Azure .....	446
18 ■ Airflow в GCP .....	465

# Содержание

---

<i>Предисловие .....</i>	14
<i>Благодарности .....</i>	16
<i>О книге .....</i>	18
<i>Об авторах .....</i>	23
<i>Об иллюстрации на обложке .....</i>	24
<b>Часть I ПРИСТУПАЕМ К РАБОТЕ .....</b>	<b>25</b>
<b>1 Знакомство с Apache Airflow .....</b>	<b>27</b>
1.1 Знакомство с конвейерами обработки данных .....	28
1.1.1 Конвейеры обработки данных как графы .....	29
1.1.2 Выполнение графа конвейера .....	30
1.1.3 Графы конвейеров и последовательные сценарии .....	32
1.1.4 Запуск конвейера с помощью диспетчеров рабочих процессов .....	33
1.2 Представляем Airflow .....	35
1.2.1 Определение конвейеров в коде ( <i>Python</i> ) гибким образом .....	35
1.2.2 Планирование и выполнение конвейеров .....	36
1.2.3 Мониторинг и обработка сбоев .....	39
1.2.4 Инкрементальная загрузка и обратное заполнение .....	41
1.3 Когда использовать Airflow .....	42
1.3.1 Причины выбрать Airflow .....	42
1.3.2 Причины не выбирать Airflow .....	43
1.4 Остальная часть книги .....	44
Резюме .....	44
<b>2 Анатомия ОАГ .....</b>	<b>46</b>
2.1 Сбор данных из множества источников .....	46
2.1.1 Изучение данных .....	47
2.2 Пишем наш первый ОАГ .....	48
2.2.1 Задачи и операторы .....	52
2.2.2 Запуск произвольного кода на Python .....	53

2.3	Запуск ОАГ в Airflow .....	56
2.3.1	Запуск Airflow в окружении Python.....	56
2.3.2	Запуск Airflow в контейнерах Docker .....	57
2.3.3	Изучаем пользовательский интерфейс Airflow .....	58
2.4	Запуск через равные промежутки времени .....	62
2.5	Обработка неудачных задач .....	64
	Резюме .....	66
<b>3</b>	<b>Планирование в Airflow .....</b>	<b>67</b>
3.1	Пример: обработка пользовательских событий .....	68
3.2	Запуск через равные промежутки времени .....	69
3.2.1	Определение интервалов .....	70
3.2.2	Интервалы на основе Cron .....	71
3.2.3	Частотные интервалы .....	73
3.3	Инкрементная обработка данных .....	74
3.3.1	Инкрементное извлечение событий.....	74
3.3.2	Динамическая привязка ко времени с использованием дат выполнения .....	75
3.3.3	Разделение данных .....	77
3.4	Даты выполнения .....	80
3.4.1	Выполнение работы с фиксированными интервалами.....	80
3.5	Использование обратного заполнения .....	82
3.5.1	Назад в прошлое .....	82
3.6	Лучшие практики для проектирования задач .....	84
3.6.1	Атомарность .....	84
3.6.2	Идемпотентность .....	86
	Резюме .....	87
<b>4</b>	<b>Создание шаблонов задач с использованием контекста Airflow .....</b>	<b>89</b>
4.1	Проверка данных для обработки с помощью Airflow.....	90
4.1.1	Определение способа загрузки инкрементальных данных .....	90
4.2	Контекст задачи и шаблонизатор Jinja .....	92
4.2.1	Создание шаблонов аргументов оператора .....	92
4.2.2	Что доступно для создания шаблонов? .....	95
4.2.3	Создание шаблона для PythonOperator .....	97
4.2.4	Предоставление переменных PythonOperator .....	102
4.2.5	Изучение шаблонных аргументов .....	104
4.3	Подключение других систем .....	105
	Резюме .....	113
<b>5</b>	<b>Определение зависимостей между задачами .....</b>	<b>114</b>
5.1	Базовые зависимости .....	115
5.1.1	Линейные зависимости .....	115
5.1.2	Зависимости «один-ко-многим» и «многие-к-одному» .....	116
5.2	Ветвление .....	119
5.2.1	Ветвление внутри задач .....	119

5.2.2	<i>Ветвление внутри ОАГ</i>	121
5.3	<b>Условные задачи</b>	126
5.3.1	<i>Условия в задачах</i>	126
5.3.2	<i>Делаем задачи условными</i>	127
5.3.3	<i>Использование встроенных операторов</i>	129
5.4	<b>Подробнее о правилах триггеров</b>	130
5.4.1	<i>Что такое правило триггеров?</i>	130
5.4.2	<i>Эффект неудач</i>	131
5.4.3	<i>Другие правила</i>	132
5.5	<b>Обмен данными между задачами</b>	133
5.5.1	<i>Обмен данными с помощью XCom</i>	134
5.5.2	<i>Когда (не) стоит использовать XCom</i>	137
5.5.3	<i>Использование настраиваемых XCom-бэкендов</i>	137
5.6	<b>Связывание задач Python с помощью Taskflow API</b>	138
5.6.1	<i>Упрощение задач Python с помощью Taskflow API</i>	139
5.6.2	<i>Когда (не) стоит использовать Taskflow API</i>	141
	<b>Резюме</b>	143

## Часть II ЗА ПРЕДЕЛАМИ ОСНОВ ..... 144

<b>6</b>	<b>Запуск рабочих процессов</b>	146
6.1	<b>Опрос условий с использованием сенсоров</b>	147
6.1.1	<i>Опрос пользовательских условий</i>	150
6.1.2	<i>Использование сенсоров в случае сбоя</i>	152
6.2	<b>Запуск других ОАГ</b>	155
6.2.1	<i>Обратное заполнение с помощью оператора TriggerDagRunOperator</i>	159
6.2.2	<i>Опрос состояния других ОАГ</i>	159
6.3	<b>Запуск рабочих процессов с помощью REST API и интерфейса командной строки</b>	163
	<b>Резюме</b>	165

<b>7</b>	<b>Обмен данными с внешними системами</b>	166
7.1	<b>Подключение к облачным сервисам</b>	167
7.1.1	<i>Установка дополнительных зависимостей</i>	168
7.1.2	<i>Разработка модели машинного обучения</i>	169
7.1.3	<i>Локальная разработка с использованием внешних систем</i>	174
7.2	<b>Перенос данных из одной системы в другую</b>	182
7.2.1	<i>Реализация оператора PostgresToS3Operator</i>	184
7.2.2	<i>Привлекаем дополнительные ресурсы для тяжелой работы</i>	187
	<b>Резюме</b>	189

<b>8</b>	<b>Создание пользовательских компонентов</b>	190
8.1	<b>Начнем с PythonOperator</b>	191
8.1.1	<i>Имитация API для рейтинга фильмов</i>	191
8.1.2	<i>Получение оценок из API</i>	194
8.1.3	<i>Создание фактического ОАГ</i>	197

<b>8.2</b>	<b>Создание собственного хука</b> .....	199
8.2.1	<i>Создание собственного хука</i> .....	200
8.2.2	<i>Создание ОАГ с помощью MovieLensHook</i> .....	206
<b>8.3</b>	<b>Создание собственного оператора</b> .....	208
8.3.1	<i>Определение собственного оператора</i> .....	208
8.3.2	<i>Создание оператора для извлечения рейтингов</i> .....	210
<b>8.4</b>	<b>Создание нестандартных сенсоров</b> .....	213
<b>8.5</b>	<b>Упаковка компонентов</b> .....	216
8.5.1	<i>Создание пакета Python</i> .....	217
8.5.2	<i>Установка пакета</i> .....	219
	<b>Резюме</b> .....	220

**9**

<b>Тестирование</b> .....	222	
<b>9.1</b>	<b>Приступаем к тестированию</b> .....	223
9.1.1	<i>Тест на благонадежность ОАГ</i> .....	223
9.1.2	<i>Настройка конвейера непрерывной интеграции и доставки</i> .....	230
9.1.3	<i>Пишем модульные тесты</i> .....	232
9.1.4	<i>Структура проекта Pytest</i> .....	233
9.1.5	<i>Тестирование с файлами на диске</i> .....	238
<b>9.2</b>	<b>Работа с ОАГ и контекстом задачи в тестах</b> .....	241
9.2.1	<i>Работа с внешними системами</i> .....	246
<b>9.3</b>	<b>Использование тестов для разработки</b> .....	254
9.3.1	<i>Тестирование полных ОАГ</i> .....	257
<b>9.4</b>	<b>Эмулируйте промышленное окружение с помощью Whirl</b> .....	257
<b>9.5</b>	<b>Создание окружений</b> .....	258
	<b>Резюме</b> .....	258

**10**

<b>Запуск задач в контейнерах</b> .....	259	
<b>10.1</b>	<b>Проблемы, вызываемые множеством разных операторов</b> .....	260
10.1.1	<i>Интерфейсы и реализации операторов</i> .....	260
10.1.2	<i>Сложные и конфликтующие зависимости</i> .....	261
10.1.3	<i>Переход к универсальному оператору</i> .....	261
<b>10.2</b>	<b>Представляем контейнеры</b> .....	262
10.2.1	<i>Что такое контейнеры?</i> .....	263
10.2.2	<i>Запуск нашего первого контейнера Docker</i> .....	264
10.2.3	<i>Создание образа Docker</i> .....	265
10.2.4	<i>Сохранение данных с использованием томов</i> .....	267
<b>10.3</b>	<b>Контейнеры и Airflow</b> .....	270
10.3.1	<i>Задачи в контейнерах</i> .....	270
10.3.2	<i>Зачем использовать контейнеры?</i> .....	270
<b>10.4</b>	<b>Запуск задач в Docker</b> .....	272
10.4.1	<i>Знакомство с DockerOperator</i> .....	272
10.4.2	<i>Создание образов для задач</i> .....	274
10.4.3	<i>Создание ОАГ с задачами Docker</i> .....	277
10.4.4	<i>Рабочий процесс на базе Docker</i> .....	280
<b>10.5</b>	<b>Запуск задач в Kubernetes</b> .....	281
10.5.1	<i>Представляем Kubernetes</i> .....	282
10.5.2	<i>Настройка Kubernetes</i> .....	283
10.5.3	<i>Использование KubernetesPodOperator</i> .....	286

10.5.4 Диагностика проблем, связанных с Kubernetes .....	290
10.5.5 Отличия от рабочих процессов на базе Docker.....	292
<b>Резюме .....</b>	<b>293</b>
<b>Часть III AIRFLOW НА ПРАКТИКЕ .....</b>	<b>294</b>
<b>11 Лучшие практики .....</b>	<b>295</b>
11.1 Написание чистых ОАГ .....	296
11.1.1 Используйте соглашения о стилях .....	296
11.1.2 Централизованное управление учетными данными .....	300
11.1.3 Единообразно указывайте детали конфигурации.....	301
11.1.4 Избегайте вычислений в определении ОАГ .....	304
11.1.5 Используйте фабричные функции для генерации распространенных шаблонов .....	306
11.1.6 Группируйте связанные задачи с помощью групп задач.....	310
11.1.7 Создавайте новые ОАГ для больших изменений .....	312
11.2 Проектирование воспроизводимых задач .....	312
11.2.1 Всегда требуйте, чтобы задачи были идемпотентными .....	312
11.2.2 Результаты задачи должны быть детерминированными .....	313
11.2.3 Проектируйте задачи с использованием парадигмы функционального программирования .....	313
11.3 Эффективная обработка данных .....	314
11.3.1 Ограничьте объем обрабатываемых данных .....	314
11.3.2 Инкрементальная загрузка и обработка .....	316
11.3.3 Кешируйте промежуточные данные .....	317
11.3.4 Не храните данные в локальных файловых системах .....	318
11.3.5 Переложите работу на внешние/исходные системы .....	318
11.4 Управление ресурсами .....	319
11.4.1 Управление параллелизмом с помощью пулов .....	319
11.4.2 Обнаружение задач с длительным временем выполнения с помощью соглашений об уровне предоставления услуг и оповещений .....	321
<b>Резюме .....</b>	<b>322</b>
<b>12 Эксплуатация Airflow в промышленном окружении .....</b>	<b>324</b>
12.1 Архитектура Airflow .....	325
12.1.1 Какой исполнитель мне подходит? .....	327
12.1.2 Настройка базы метаданных для Airflow .....	328
12.1.3 Присмотримся к планировщику .....	330
12.2 Установка исполнителей .....	334
12.2.1 Настройка SequentialExecutor .....	335
12.2.2 Настройка LocalExecutor .....	335
12.2.3 Настройка CeleryExecutor .....	336
12.2.4 Настройка KubernetesExecutor .....	339
12.3 Работа с журналами всех процессов Airflow .....	347
12.3.1 Вывод веб-сервера .....	347
12.3.2 Вывод планировщика .....	348

12.3.3	<i>Журналы задач</i> .....	349
12.3.4	<i>Отправка журналов в удаленное хранилище</i> .....	350
12.4	<b>Визуализация и мониторинг метрик Airflow</b> .....	350
12.4.1	<i>Сбор метрик из Airflow</i> .....	351
12.4.2	<i>Настройка Airflow для отправки метрик</i> .....	353
12.4.3	<i>Настройка Prometheus для сбора метрик</i> .....	353
12.4.4	<i>Создание дашбордов с Grafana</i> .....	356
12.4.5	<i>Что следует мониторить?</i> .....	358
12.5	<b>Как получить уведомление о невыполненной задаче</b> .....	360
12.5.1	<i>Оповещения в ОАГ и операторах</i> .....	360
12.5.2	<i>Определение соглашений об уровне предоставления услуги</i> .....	362
12.6	<b>Масштабируемость и производительность</b> .....	364
12.6.1	<i>Контроль максимального количества запущенных задач</i> .....	365
12.6.2	<i>Конфигурации производительности системы</i> .....	366
12.6.3	<i>Запуск нескольких планировщиков</i> .....	367
	<b>Резюме</b> .....	368
<b>13</b>	<b>Безопасность в Airflow</b> .....	369
13.1	<b>Обеспечение безопасности веб-интерфейса Airflow</b> .....	370
13.1.1	<i>Добавление пользователей в интерфейс RBAC</i> .....	371
13.1.2	<i>Настройка интерфейса RBAC</i> .....	374
13.2	<b>Шифрование хранимых данных</b> .....	375
13.2.1	<i>Создание ключа Fernet</i> .....	375
13.3	<b>Подключение к службе LDAP</b> .....	377
13.3.1	<i>Разбираемся с LDAP</i> .....	378
13.3.2	<i>Извлечение пользователей из службы LDAP</i> .....	380
13.4	<b>Шифрование трафика на веб-сервер</b> .....	381
13.4.1	<i>Разбираемся с протоколом HTTP</i> .....	381
13.4.2	<i>Настройка сертификата для HTTPS</i> .....	384
13.5	<b>Извлечение учетных данных из систем управления секретами</b> .....	388
	<b>Резюме</b> .....	392
<b>14</b>	<b>Проект: поиск самого быстрого способа передвижения по Нью-Йорку</b> .....	393
14.1	<b>Разбираемся с данными</b> .....	396
14.1.1	<i>Файловый ресурс Yellow Cab</i> .....	397
14.1.2	<i>REST API Citi Bike</i> .....	397
14.1.3	<i>Выбор плана подхода</i> .....	399
14.2	<b>Извлечение данных</b> .....	400
14.2.1	<i>Скачиваем данные по Citi Bike</i> .....	400
14.2.2	<i>Загрузка данных по Yellow Cab</i> .....	402
14.3	<b>Применение аналогичных преобразований к данным</b> .....	405
14.4	<b>Структурирование конвейера обработки данных</b> .....	410
14.5	<b>Разработка идемпотентных конвейеров обработки данных</b> .....	411
	<b>Резюме</b> .....	414

**Часть IV ОБЛАКО .....** 415

**15 Airflow и облако .....** 417

15.1	Проектирование стратегий (облачного) развертывания .....	418
15.2	Операторы и хуки, предназначенные для облака.....	420
15.3	Управляемые сервисы.....	421
	15.3.1 <i>Astronomer.io</i> .....	421
	15.3.2 <i>Google Cloud Composer</i> .....	422
	15.3.3 <i>Amazon Managed Workflows for Apache Airflow</i> .....	423
15.4	Выбор стратегии развертывания .....	423
	Резюме .....	425

**16 Airflow и AWS .....** 426

16.1	Развертывание Airflow в AWS .....	426
	16.1.1 Выбор облачных сервисов .....	427
	16.1.2 Проектирование сети.....	428
	16.1.3 Добавление синхронизации ОАГ.....	430
	16.1.4 Масштабирование с помощью CeleryExecutor .....	430
	16.1.5 Дальнейшие шаги .....	432
16.2	Хуки и операторы, предназначенные для AWS .....	432
16.3	Пример использования: бессерверное ранжирование фильмов с AWS Athena .....	434
	16.3.1 Обзор .....	434
	16.3.2 Настройка ресурсов .....	435
	16.3.3 Создание ОАГ.....	438
	16.3.4 Очистка .....	445
	Резюме .....	445

**17 Airflow и Azure .....** 446

17.1	Развертывание Airflow в Azure .....	446
	17.1.1 Выбор сервисов .....	447
	17.1.2 Проектирование сети.....	448
	17.1.3 Масштабирование с помощью CeleryExecutor .....	449
	17.1.4 Дальнейшие шаги .....	450
17.2	Хуки и операторы, предназначенные для Azure .....	451
17.3	Пример: бессерверное ранжирование фильмов с Azure Synapse .....	452
	17.3.1 Обзор .....	452
	17.3.2 Настройка ресурсов .....	453
	17.3.3 Создание ОАГ.....	457
	17.3.4 Очистка .....	463
	Резюме .....	464

**18 Airflow в GCP .....** 465

18.1	Развертывание Airflow в GCP .....	465
	18.1.1 Выбор сервисов .....	466
	18.1.2 Развертывание в GKE с помощью Helm .....	468

18.1.3 Интеграция с сервисами Google.....	471
18.1.4 Проектирование сети.....	472
18.1.5 Масштабирование с помощью CeleryExecutor .....	473
18.2 Хуки и операторы, предназначенные для GCP .....	476
18.3 Пример использования: бессерверный рейтинг фильмов в GCP.....	481
18.3.1 Загрузка в GCS.....	481
18.3.2 Загрузка данных в BigQuery.....	483
18.3.3 Извлечение рейтингов, находящихся в топе .....	485
Резюме .....	488
Приложение A Запуск примеров кода .....	490
Приложение B Структуры пакетов Airflow 1 и 2 .....	494
Приложение C Сопоставление метрик в Prometheus .....	498
Предметный указатель .....	500